

A Theoretical Analysis of NDCG Type Ranking Measures

Yining Wang (antoniowyn@gmail.com)

Institute for Interdisciplinary Information Sciences, Tsinghua University,
Beijing, P.R.China

Liwei Wang (wanglw@cis.pku.edu.cn)

School of Electronics Engineering and Computer Science,
Peking University
Beijing, P.R.China

Yuanzhi Li (invinciblec.lee@gmail.com)

Institute for Interdisciplinary Information Sciences, Tsinghua University,
Beijing, P.R.China

Di He (wolfink@gmail.com)

School of Electronics Engineering and Computer Science,
Peking University
Beijing, P.R.China

Tie-Yan Liu (Tie-Yan.Liu@microsoft.com)

Microsoft Research Asia,
Beijing, P.R. China

Wei Chen (wche@microsoft.com)

Microsoft Research Asia,
Beijing, P.R. China

April 25, 2013

Abstract

A central problem in ranking is to design a ranking measure for evaluation of ranking functions. In this paper we study, from a theoretical perspective, the widely used Normalized Discounted Cumulative Gain (NDCG)-type ranking measures. Although there are extensive empirical studies of NDCG, little is known about its theoretical properties. We first show that, whatever the ranking function is, the standard NDCG which adopts a logarithmic discount, converges to 1 as the number of items to rank goes to infinity. On the first sight, this result is very surprising. It seems to imply that NDCG cannot differentiate good and bad ranking functions, contradicting to the empirical success of NDCG in many applications. In order to have a deeper understanding of ranking measures in general, we propose a notion referred to as *consistent distinguishability*. This notion captures the intuition that a ranking measure should have such a property: For every pair of substantially different ranking functions, the ranking measure can decide which one is better in a *consistent* manner on almost all datasets. We show that NDCG with logarithmic discount has consistent distinguishability although it converges to the same limit for all ranking functions. We next characterize the set of all feasible discount functions for NDCG according to the concept of consistent distinguishability. Specifically we show that whether NDCG has consistent distinguishability depends on how fast the discount decays, and r^{-1} is a critical point. We then turn to the cut-off version of NDCG, i.e., NDCG@k. We analyze the distinguishability of NDCG@k for various choices of k and the discount functions. Experimental results on real Web search datasets agree well with the theory.

1 Introduction

Ranking has been extensively studied in information retrieval, machine learning and statistics. It plays a central role in various applications such as search engine, recommendation system, expert finding, to name a few. In many situations one wants to have, by learning, a good ranking function [15, 18, 23]. Thus a fundamental problem is how to design a ranking measure to evaluate the performance of a ranking function.

Unlike classification and regression for which there are simple and natural performance measures, evaluating ranking functions has proved to be more difficult. Suppose there are n objects to rank. A ranking evaluation measure must induce a total order on the $n!$ possible ranking results. There seem to be many ways to define ranking measures and several evaluation measures have been proposed [11, 32, 4, 1, 28]. In fact, as pointed out by some authors, there is no single optimal ranking measure that works for any application [16].

The focus of this work is the Normalized Discounted Cumulative Gain (NDCG) which is one of the most popular evaluation measures in Web search [21, 22]. NDCG has two advantages compared to many other measures. First, NDCG allows each retrieved document has graded relevance while most traditional ranking measures only allow binary relevance. That is, each document is viewed as either relevant or not relevant by previous ranking measures, while there can be degrees of relevancy for documents in NDCG. Second, NDCG involves a discount function over the rank while many other measures uniformly weight all positions. This feature is particularly important for search engines as users care top ranked documents much more than others.

The importance of NDCG as well as other ranking measures in modern search engines is not limited as evaluation metrics. Currently ranking measures are also used as guidance for design of ranking functions due to works from the learning to rank area. Although early results of learning to rank often reduce ranking problem to classification or regression [15, 18, 23, 26, 5], recently there is evidence that learning a ranking function by optimizing a ranking measure such as NDCG is a promising approach [33, 36]. However, using the ranking measure as objective function to optimize is computationally intractable. Inspired by approaches in classification, some state of the art algorithms optimize a surrogate loss instead [8, 35].

In the past a few years, there is rapidly growing interest in studying consistency of learning to rank algorithms that optimize surrogate losses. Such studies are motivated by the research of consistency of surrogate losses for classification [38, 6, 37, 31], which is a well-established theory in machine learning. Consistency of ranking is more complicated than classification as there are more than one possible ranking measures. One needs to study consistency with respect to a specific ranking measure. That is, whether the minimization of the surrogate leads to optimal predictions according to the risk defined by the given evaluation measure.

The research of consistency for ranking was initiated in [14, 17]. In fact, [17] showed that no convex surrogate loss can be consistent with the Pairwise Disagreement (PD) measure. This result was further generalized in [7, 9], where non-existence of convex surrogate loss with Average Precision and Expected Reciprocal Rank were proved.

In contrast to the above negative results, [27] showed that there do exist NDCG consistent surrogates. Furthermore, by using a slightly stronger notion of NDCG consistency they showed that any NDCG consistent surrogate must be a Bregman distance. In a sense, these results mean that NDCG is a good ranking measure from a learning-to-rank point of view.

NDCG is a normalization of the Discounted Cumulative Gain (DCG) measure. (For formal definition of both DCG and NDCG, please see Section 2.) DCG is a weighted sum of the degree of relevancy of the ranked items. The weight is a decreasing function of the rank (position) of the object, and therefore called discount. The original reason for introducing the discount is that the probability that a user views a document decreases with respect to its rank. NDCG normalizes DCG by the Ideal DCG (IDCG), which is simply the DCG measure of the best ranking result. Thus NDCG measure is always a number in $[0, 1]$. Strictly speaking, NDCG is a family of ranking

measures, since there is flexibility in choosing the discount function. The logarithmic discount $\frac{1}{\log(1+r)}$, where r is the rank, dominated the literature and applications. We will refer to NDCG with logarithmic discount as the *standard* NDCG. Another discount function appeared in literature is r^{-1} , which is called Zipfian in Information Retrieval [24]. Search engine systems also use a cut-off top-k version of NDCG. That is, the discount is set to be zero for ranks larger than k . Such NDCG measure is usually referred to as NDCG@k.

Given the importance and popularity of NDCG, there have been extensive studies on this measure, mainly in the field of Information Retrieval [2, 24, 3, 34, 29]. All these research are conducted from an empirical perspective by doing experiments on benchmark datasets. Although these works gained insights about NDCG, there are still important issues unaddressed. We list a few questions that naturally arise.

- As pointed out in [16], there has not been any theoretically sound justification for using a logarithmic ($\frac{1}{\log(1+r)}$) discount other than the fact that it is a smooth decay.
- Is it possible to characterize the class of discount functions that are feasible for NDCG?
- For the standard NDCG@k, the discount is a combination of a very slow logarithmic decay and a hard cut-off. Why don't simply use a smooth discount that decays fast?

In this paper, we study the NDCG type ranking measures and address the above questions from a theoretical perspective. The goal of our study is twofold. First, we aim to provide a better understanding and theoretical justification of NDCG as an evaluation measure. Second, we hope that our results would shed light and be useful for further research on learning to rank based on NDCG. Specifically we analyze the behavior of NDCG as the number of objects to rank getting large. Asymptotics, including convergence and asymptotic normality, of many traditional ranking measures have been studied in depth in statistics, especially for Linear Rank Statistics and measures that are U-statistics [19, 25]. [12] observed that ranking measures such as Area under the ROC Curve (AUC), P-Norm Push and DCG can be viewed as Conditional Linear Rank Statistics. That is, conditioned on the relevance degrees of the items, these measures are Linear Rank Statistics [19]. They show uniform convergence based on an orthogonal decomposition of the measure. The convergence relies on the fact that the measure can be represented as a (conditional) average of a fixed score-generating function. Part of our work consider the convergence of NDCG and are closely related to [12]. However, their results do not apply to our problem, because the score-generating function for NDCG is not fixed, it changes with the number of objects.

1.1 Our Results

Our study starts from an analysis of the standard NDCG (i.e., the one using logarithmic discount). The first discovery is that for *every* ranking function, the NDCG measure converges to 1 as the number of items to rank goes to infinity. This result is surprising. On the first sight it seems to mean that the widely used standard NDCG cannot differentiate good and bad ranking systems when the data is of large size. This problem may be serious because huge dataset is common in applications such as Web search.

To have a deeper understanding of NDCG, we first study what are the desired properties a good ranking measure should have. In this paper we propose a notion referred to as *consistent distinguishability*, which we believe that every ranking measure needs to have. Before describing the definition of consistent distinguishability, let us see a motivating example. Suppose we want to select, from two ranking functions f_1, f_2 , a better one on ranking “sea” images (that is, if an image contains sea, we hope it is ranked near the top). Since there are billions of sea images on the web, a commonly used method is to randomly draw, say, a million data and evaluate the two functions on them. A crucial assumption underlying this approach is that the evaluation result will be “stable”

on large datasets. That is, if on this randomly drawn dataset f_1 is better than f_2 according to the ranking measure, then with high probability over the random draw of another large dataset, f_1 should still be better than f_2 . In other words, f_1 is *consistently* better than f_2 according to the ranking measure.

Our definition of consistent distinguishability captures the above intuition. It requires that for two substantially different ranking functions, the ranking measure can decide which one is better consistently on almost all datasets. (See Definition 3 for formal description.) In a broader sense, consistent distinguishability is a desired property to all performance statistics (not only to ranking). For classification and regression, this property trivially holds because of the simplicity of the evaluation measures. For ranking however, things are much more complicated. It is not a priori clear whether important ranking measures such as NDCG have consistent distinguishability.

Our next main result shows that although the standard NDCG always converges to 1, it can consistently distinguish every pair of substantially different ranking functions. Therefore, if one ignores the numerical scaling problem, standard NDCG is a good ranking measure.

We then study NDCG with other possible discount. We characterize the class of discount functions that are feasible for NDCG. It turns out that the Zipfian r^{-1} is a critical point. If a discount function decays slower than r^{-1} , the resulting NDCG measure has strong power of consistent distinguishability. If a discount decays substantially faster than r^{-1} , then it does not have this desired property. Even more, such ranking measures do not converge as the number of objects to rank goes to infinity.

Interestingly, this characterization result also provides a better understanding of the cut-off version NDCG@k. In particular, it gives a theoretical explanation to the previous question that why popular NDCG@k uses a combination of slow logarithmic decay and a hard cut-off as its discount rather than a smooth discount which decays fast.

Finally we consider how to choose the cut-off threshold for NDCG@k from the distinguishability point of view. We analyze the behavior of the measure for various choices of k as well as the discount. We suggest that choosing k as certain function of the size of the dataset may be appropriate.

The rest of this paper is organized as follows. Section 2 provides basic notions and definitions. Section 3 contains the main theorems and key lemmas for the distinguishability theorem. The experimental results are given in 4. All proofs are given in Appendix A-E.

2 Preliminaries

Let \mathcal{X} be the instance space, and let x_1, \dots, x_n ($x_i \in \mathcal{X}$) be n objects to rank. Let \mathcal{Y} be a finite set of degrees of relevancy. The simplest case is $\mathcal{Y} = \{0, 1\}$, where 0 corresponds to “irrelevant” and 1 corresponds to “relevant”. Generally \mathcal{Y} may contain more numbers; and for $y \in \mathcal{Y}$, the larger y is, the more relevant it represents. Let f be a ranking function¹. We assume that f is a mapping from \mathcal{X} to \mathbb{R} . For each object $x \in \mathcal{X}$, f gives it a score $f(x)$. For n objects x_1, \dots, x_n , f ranks them according to their scores $f(x_1), \dots, f(x_n)$. The resulting ranking list, denoted by $x_{(1)}^f, \dots, x_{(n)}^f$, satisfies $f(x_{(1)}^f) \geq \dots \geq f(x_{(n)}^f)$.

Let y_1, \dots, y_n ($y_i \in \mathcal{Y}$) be the degree of relevancy associated with x_1, \dots, x_n . We will denote by $S_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ the set of data to rank. As in existing literature [18, 13], we assume that $(x_1, y_1), \dots, (x_n, y_n)$ are i.i.d. sample drawn from an underlying distribution P_{XY} over $\mathcal{X} \times \mathcal{Y}$. Also let $y_{(1)}^f, \dots, y_{(n)}^f$ be the corresponding relevancy of $x_{(1)}^f, \dots, x_{(n)}^f$.

¹The ranking function we defined is often called scoring function in literature; and ranking function has a more general definition: For fixed n , a general ranking function can be any permutation on $[n]$. However, scoring functions are used by most search engines. Also in this paper we study the behavior of the ranking measure of a fixed ranking function as n grows, so we focus on scoring functions. But note that Theorem 1 and Theorem 6 hold for any sequence of general ranking functions.

The following is the formal definition of NDCG. Here we give a slightly simplified version tailored to our problem.

Definition 1. Let $D(r)$ ($r \geq 1$) be a discount function. Let f be a ranking function, and S_n be a dataset. The Discounted Cumulative Gain (DCG) of f on S_n with discount D is defined as²

$$\text{DCG}_D(f, S_n) = \sum_{r=1}^n y_{(r)}^f D(r). \quad (1)$$

Let the Ideal DCG defined as $\text{IDCG}_D(S_n) = \max_{f'} \sum_{r=1}^n y_{(r)}^{f'} D(r)$ be the DCG value of the best ranking function on S_n .

The NDCG of f on S_n with discount D is defined as

$$\text{NDCG}_D(f, S_n) = \frac{\text{DCG}_D(f, S_n)}{\text{IDCG}_D(S_n)}. \quad (2)$$

We call NDCG *standard*, if its associated discount function is the inverse logarithm decay $D(r) = \frac{1}{\log(1+r)}$. Note that the base of the logarithm does not matter for NDCG, since constant scaling will cancel out due to normalization. We will assume it is the natural logarithm throughout this paper.

An important property of eq.(2) is that if a ranking function f' preserves the order of the ranking function f , then $\text{NDCG}_D(f', S_n) = \text{NDCG}_D(f, S_n)$ for all S_n . Here by preserving order we mean that for $\forall x, x' \in \mathcal{X}$, $f(x) > f(x')$ implies $f'(x) > f'(x')$, and vice versa. Thus the ranking measure NDCG is not just defined on a single function f , but indeed defined on an equivalent class of ranking functions which preserve order of each other.

Below we will frequently use a special ranking function \tilde{f} that preserves the order of f .

Definition 2. Let f be a ranking function. We call \tilde{f} the canonical version of f , which is defined as

$$\tilde{f}(x) = \Pr_{X \sim P_X} [f(X) \leq f(x)].$$

The canonical \tilde{f} has the following properties, which can be easily proved by the definition.

Lemma 1. For every ranking function f , its canonical version \tilde{f} preserves the order of f . In addition, $\tilde{f}(X)$ has uniform distribution on $[0, 1]$.

Finally, we point out that although originally the discount $D(r)$ is defined on positive integers r , below we will often treat $D(r)$ as a function of a real variable. That is, we view r take nonnegative real values. We will also consider derivative and integral of $D(r)$, denoted by $D'(r)$ and $\int D(r)dr$ respectively.

3 Main Results

In this section, we give the main results of the paper. In Section 3.1 we study the standard NDCG, i.e., NDCG with logarithmic discount. In Section 3.2 we consider feasible discount other than the standard logarithmic one. We analyze the top-k cut-off version NDCG@k in Section 3.3. For clarity reasons, some of the results in Section 3.1, 3.2, and 3.3 are given for the simplest case that the relevance score is binary. Section 3.4 provides complete results for the general case.

²Usually DCG is defined as $\text{DCG}_D(f, S_n) = \sum_{r=1}^n G(y_{(r)}^f) D(r)$, where G is a monotone increasing function (e.g., $G(y) = 2^y - 1$). Here we omit G for notational simplicity. This does not lose any generality as we can assume that \mathcal{Y} changes to $G(\mathcal{Y})$.

3.1 Standard NDCG

To study the behavior of the standard NDCG, we first consider the limit of this measure when the number of objects to rank goes to infinity. As stated in Section 2, we assume the data are i.i.d. drawn from some fixed underlying distribution. Surprisingly, it is easy to show that for every ranking function, standard NDCG converges to 1 almost surely.

Theorem 1. *Let $D(r) = \frac{1}{\log(1+r)}$. Then for every ranking function f ,*

$$\text{NDCG}_D(f, S_n) \rightarrow 1, \quad a.s.$$

The proof is given in Appendix D.

At the first glance, the above result is quite negative for standard NDCG. It seems to say that in the limiting case, standard NDCG cannot differentiate ranking functions. However, Theorem 1 only considers the limits. To have a better understanding of NDCG, we need to make a deeper analysis of its power of distinguishability. In particular, Theorem 1 does not rule out the possibility that the standard NDCG can consistently distinguish substantially different ranking functions. Below we give the formal definition that two ranking functions are consistently distinguishable by a ranking measure \mathcal{M} .

Definition 3. *Let $(x_1, y_1), (x_2, y_2), \dots$ be i.i.d. instance-label pairs drawn from the underlying distribution P_{XY} over $\mathcal{X} \times \mathcal{Y}$. Let $S_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$. A pair of ranking functions f_0, f_1 is said to be consistently distinguishable by a ranking measure \mathcal{M} , if there exists a negligible function³ $\text{neg}(N)$ and $b \in \{0, 1\}$ such that for every sufficiently large N , with probability $1 - \text{neg}(N)$,*

$$\mathcal{M}(f_b, S_n) > \mathcal{M}(f_{1-b}, S_n),$$

holds for all $n \geq N$ simultaneously.

Consistent distinguishability is appealing. One would like a ranking measure \mathcal{M} to have the property that every two substantially different ranking functions are consistently distinguishable by \mathcal{M} . The next theorem shows that standard NDCG does have such a desired property. For clarity, here we state the theorem for the simple binary relevance case, i.e., $\mathcal{Y} = \{0, 1\}$. It is easy to extend the result to the general case that \mathcal{Y} is any finite set.

Theorem 2. *For every pair of ranking functions f_0, f_1 , let $\bar{y}^{f_i}(s) = \Pr[Y = 1 | \tilde{f}_i(X) = s]$, $i = 0, 1$. Assume $\bar{y}^{f_0}(s)$ and $\bar{y}^{f_1}(s)$ are Hölder continuous in s . Then, unless $\bar{y}^{f_0}(s) = \bar{y}^{f_1}(s)$ almost everywhere on $[0, 1]$, f_0 and f_1 are consistently distinguishable by standard NDCG.*

The proof is given in Appendix A.

Theorem 2 provides theoretical justification for using standard NDCG as a ranking measure, and answers the first question raised in Introduction. Although standard NDCG converges to the same limit for all ranking functions, it is still a good ranking measure with strong consistent distinguishability (if we ignore the numerical scaling issue).

3.2 Characterization of Feasible Discount Functions

In the previous section we demonstrate that standard NDCG is a good ranking measure. In both literatures and real applications, standard NDCG is dominant. However, there is no known theoretical evidence that the logarithmic function is the only feasible discount, or it is the optimal one. In this subsection, we will investigate other discount functions. We study the asymptotic behavior and

³A negligible function $\text{neg}(N)$ means that for $\forall c$, $\text{neg}(N) < N^{-c}$ for sufficiently large N .

distinguishability of the induced NDCG measures and compare to the standard NDCG. Finally, we will characterize the class of discount functions which we think are feasible for NDCG. For the sake of clarity, the results in this subsection are given for the simplest case that $\mathcal{Y} = \{0, 1\}$. Complete results will be given in Section 3.4.

Standard NDCG utilizes the logarithmic discount which decays slowly. In the following we first consider a discount that decays a little faster. Specifically we consider $D(r) = r^{-\beta}$ ($0 < \beta < 1$). Let us first investigate the limit of the ranking measure as the number of objects goes to infinity.

Theorem 3. *Assume $D(r) = r^{-\beta}$ where $\beta \in (0, 1)$. Assume also $p = \Pr[Y = 1] > 0$ and $\bar{y}^f(s) = \Pr[Y = 1 | \tilde{f}(X) = s]$ is a continuous function. Then*

$$\text{NDCG}_D(f, S_n) \xrightarrow{p} \frac{(1 - \beta) \int_0^1 \bar{y}^f(s) \cdot (1 - s)^{-\beta} ds}{p^{1-\beta}}. \quad (3)$$

The proof will be given in Appendix D.

For $D(r) = r^{-\beta}$ ($\beta \in (0, 1)$), NDCG no longer converges to the same limit for all ranking functions. The limit is actually a correlation between $\bar{y}^f(s)$ and $(1 - s)^{-\beta}$. For a good ranking function f , $\bar{y}^f(s) = \Pr[Y = 1 | \tilde{f}(X) = s]$ is likely to be an increasing function of s , and thus has positive correlation with $(1 - s)^{-\beta}$. Therefore, the limit of the ranking measure already differentiate good and bad ranking functions to some extent.

We next study whether NDCG with polynomial discount has power of distinguishability as strong as the standard NDCG. That is, we will see if Theorem 2 holds for NDCG with $r^{-\beta}$ ($\beta \in (0, 1)$).

Theorem 4. *Let $D(r) = r^{-\beta}$, $\beta \in (0, 1)$. Assume $p = \Pr[Y = 1] > 0$. For every pair of ranking functions f_0, f_1 , denote $\bar{y}^{f_i}(s) = \Pr[Y = 1 | \tilde{f}_i(X) = s]$, $i = 0, 1$, and $\Delta y(s) = \bar{y}^{f_0}(s) - \bar{y}^{f_1}(s)$. Suppose at least one of the following two conditions hold: 1) $\int_0^1 \Delta y(s)(1 - s)^{-\beta} ds \neq 0$; 2) $\bar{y}^{f_0}(s)$, $\bar{y}^{f_1}(s)$ are Hölder continuous with Hölder continuity constant α satisfying $\alpha > 3(1 - \beta)$, and $\Delta y(1) \neq 0$. Then f_0 and f_1 are strictly distinguishable with high probability by NDCG with discount $D(r)$.*

The proof will be given in Appendix E.

Theorem 4 involves two conditions. Satisfying either of them leads to strictly distinguishable with high probability. The first condition simply means that $\text{NDCG}_D(f_0, S_n)$ and $\text{NDCG}_D(f_1, S_n)$ converge to different limits and therefore the two functions are consistently distinguishable in the strongest sense. The second condition deals with the case that $\text{NDCG}_D(f_0, S_n)$ and $\text{NDCG}_D(f_1, S_n)$ converge to the same limit. Comparing the distinguishability of NDCG with $r^{(-\beta)}$ discount with the standard NDCG, in most cases $r^{(-\beta)}$ discount has stronger distinguishability than standard NDCG (i.e., when the measures of two ranking functions converge to different limits). On the other hand, if we consider the worst case, standard NDCG is better, because it requires less conditions for consistent distinguishability.

We next study the Zipfian discount $D(r) = r^{-1}$. The following theorem describes the limit of the ranking measure.

Theorem 5. *Assume $D(r) = r^{-1}$. Assume also $p = \Pr[Y = 1] > 0$ and $\bar{y}^f(s) = \Pr[Y = 1 | \tilde{f}(X) = s]$ is a continuous function. Then*

$$\text{NDCG}_D(f, S_n) \xrightarrow{p} \Pr[Y = 1 | \tilde{f}(X) = 1]. \quad (4)$$

The proof of Theorem 5 will be given in Appendix D.

The limit of NDCG with Zipfian discount depends only on the performance of the ranking function for the top ranks. The relevancy of lower ranked items does not affect the limit.

The next logical step would be analyzing the power of distinguishability of NDCG with Zipfian discount. However we are not able to prove that consistent distinguishability holds for this ranking measure. The techniques developed for distinguishability theorems given above does not apply to the Zipfian discount. Although we cannot disprove it distinguishability, we suspect that Zipfian does not have strong consistent distinguishability power.

Finally, we consider discount functions that decay substantially faster than r^{-1} . We will show that with these discount, NDCG does not converge as the number of objects tends to infinity. More importantly, such NDCG does not have the desired consistent distinguishability property.

Theorem 6. *Let \mathcal{X} be instance space. For any $x \in \mathcal{X}$, let $y_x^* = \operatorname{argmax}_{y \in \mathcal{Y}} \Pr(Y = y|X = x)$. Assume that there is an absolute constant $\delta > 0$ such that for every $x \in \mathcal{X}$, $\Pr(Y = y|X = x) \geq \delta \cdot \Pr(Y = y_x^*|X = x)$ for all $y \in \mathcal{Y}$. If $\sum_{r=1}^{\infty} D(r) \leq B$ for some constant $B > 0$, then $\text{NDCG}_D(f, S_n)$ does not converge in probability for any ranking function f . In particular, if $D(r) \leq r^{-(1+\epsilon)}$ for some $\epsilon > 0$, $\text{NDCG}_D(f, S_n)$ does not converge. Moreover, every pair of ranking functions are not consistently distinguishable by NDCG with such discount.*

The proof is given in Appendix D.

Now we are able to *characterize* the feasible discounts for NDCG according to the results given so far. The logarithmic $\frac{1}{\log(1+r)}$ and polynomial $r^{-\beta}$ ($\beta \in (0, 1)$) are feasible discount functions for NDCG. For different ranking functions, standard NDCG converges to the same limit while the $r^{-\beta}$ ($\beta \in (0, 1)$) one converges to different limits in most cases. However, if we ignore the numerical scaling issue, both logarithmic and $r^{-\beta}$ ($\beta \in (0, 1)$) discount have consistent distinguishability. The Zipfian r^{-1} discount is on the borderline. It is not clear whether it has strong power of distinguishability. Discount that decays faster than $r^{-(1+\epsilon)}$ for some $\epsilon > 0$ is not appropriate for NDCG when the data size is large.

3.3 Cut-off Versions of NDCG

In this section we study the top- k version of NDCG, i.e., $\text{NDCG}@k$. For $\text{NDCG}@k$, the discount function is set as $D(r) = 0$ for all $r > k$. The motivation of using $\text{NDCG}@k$ is to pay more attention to the top-ranked results. Logarithmic discount is also dominant for $\text{NDCG}@k$. We will call this measure standard $\text{NDCG}@k$. As already stated in Introduction, a natural question of standard $\text{NDCG}@k$ is why use a combination of a very low logarithmic decay and a hard cut-off as the discount function. Why not simply use a smooth discount with fast decay, which seems more natural. In fact, this question has already been answered by Theorem 6. NDCG with such discount does not have strong power of distinguishability.

We next address the issue that how to choose the cut-off threshold k . It is obvious that setting k as a constant independent of n is not appropriate, because the partial sum of the discount is bounded and according to Theorem 6 the ranking measure does not converge. So k must grow unboundedly as n goes to infinity. Below we investigate the convergence and distinguishability of $\text{NDCG}@k$ for various choices of k and the discount function. For clarity reason we assume here $\mathcal{Y} = \{0, 1\}$, and general results will be given in Section 3.4. The proofs of all theorems in this section will be given in Appendix D. We first consider the case $k = o(n)$.

Theorem 7. *Let $\mathcal{Y} = \{0, 1\}$. Assume $D(r)$ is a discount function and $\sum_{r=1}^{\infty} D(r)$ is unbounded. Suppose $k = o(n)$ and $k \rightarrow \infty$ as $n \rightarrow \infty$. Let $\tilde{D}(r) = D(r)$ for all $r \leq k$ and $\tilde{D}(r) = 0$ for all $r > k$. Assume also that $p = \Pr[Y = 1] > 0$ and $\bar{y}^f(s) = \Pr[Y = 1|\tilde{f}(X) = s]$ is a continuous function. Then*

$$\text{NDCG}_{\tilde{D}}(f, S_n) \xrightarrow{p} \Pr[Y = 1|\tilde{f}(X) = 1]. \quad (5)$$

The limit of NDCG@k where $k = o(n)$ is exactly the same as NDCG with Zipfian discount. Also like the Zipfian, the distinguishability power of this NDCG@k measure is not clear.

We next consider the case $k = cn$ for some constant $c \in (0, 1)$. We study the standard logarithmic and the polynomial discount respectively in the following two theorems.

Theorem 8. Assume $D(r) = \frac{1}{\log(1+r)}$ and $\mathcal{Y} = \{0, 1\}$. Let $k = cn$ for some constant $c \in (0, 1)$. Define the cut-off discount function \tilde{D} as $\tilde{D}(r) = D(r)$ if $r \leq k$ and $\tilde{D}(r) = 0$ otherwise. Assume also $p = \Pr[Y = 1] > 0$ and $\bar{y}^f(s) = \Pr[Y = 1 | \tilde{f}(X) = s]$ is a continuous function. Then

$$\text{NDCG}_{\tilde{D}}(f, S_n) \xrightarrow{p} \frac{c}{\min\{c, p\}} \cdot \Pr[Y = 1 | \tilde{f}(X) \geq 1 - c]. \quad (6)$$

Theorem 9. Assume $D(r) = r^{-\beta}$ and $\mathcal{Y} = \{0, 1\}$, where $\beta \in (0, 1)$. Let $k = cn$ for some constant $c \in (0, 1)$. Define the cut-off discount function $\tilde{D}(r) = D(r)$ if $r \leq k$ and $\tilde{D}(r) = 0$ otherwise. Assume also $p = \Pr[Y = 1] > 0$ and $\bar{y}^f(s) = \Pr[Y = 1 | r(X) = s]$ is a continuous function. Then

$$\text{NDCG}_{\tilde{D}}(f, S_n) \xrightarrow{p} \frac{1 - \beta}{(\min\{c, p\})^{1-\beta}} \cdot \int_{1-c}^1 \bar{y}^f(s) \cdot (1 - s)^{-\beta} ds. \quad (7)$$

The consistent distinguishability of the two measures considered in Theorem 8 and Theorem 9 are similar to their corresponding full NDCG respectively. To be precise, for NDCG@k ($k = cn$) with logarithmic discount and NDCG@k with $r^{-\beta}$ ($\beta \in (0, 1)$) discount, consistent distinguishability holds under the condition given in Theorem 2 and Theorem 4 respectively. Hence these two cut-off versions NDCG are feasible ranking measures.

3.4 Results for General \mathcal{Y}

Some theorems given so far assume $\mathcal{Y} = \{0, 1\}$. Here we give complete results for the general case that $|\mathcal{Y}| \geq 2$, and $\mathcal{Y} = \{\eta_1, \dots, \eta_{|\mathcal{Y}|}\}$. We only state the theorems and omit the proofs, which are straightforward modifications of the special case $\mathcal{Y} = \{0, 1\}$. The case $D(r) = \frac{1}{\log(1+r)}$ has already been included in Theorem 1. It always converges to 1 whatever the ranking function is. We next consider $r^{-\beta}$ decay.

Theorem 10. Assume $D(r) = r^{-\beta}$ with $\beta \in (0, 1)$. Suppose that $\mathcal{Y} = \{\eta_1, \dots, \eta_{|\mathcal{Y}|}\}$, where $\eta_1 > \dots > \eta_{|\mathcal{Y}|}$. Assume $f(X) \in [a, b]$; $f(X)$ has a probability density function such that $\mathbb{P}(f(X) = s) > 0$ for all $s \in [a, b]$; $\Pr(Y = \eta_j) > 0$ and $\Pr(Y = \eta_j | \tilde{f}(X) = s)$ is a continuous function of s for all j . Then

$$\text{NDCG}_D(f, S_n) \xrightarrow{p} \frac{(1 - \beta) \int_0^1 \mathbb{E}[Y | \tilde{f}(X) = s] (1 - s)^{-\beta} ds}{\sum_{j=1}^{|\mathcal{Y}|} \eta_j (R_j^{1-\beta} - R_{j-1}^{1-\beta})}$$

where $R_0 = 0$; $R_j = \Pr(Y \geq \eta_j)$.

The next theorem is for top-k type NDCG measures, where $k = o(n)$.

Theorem 11. Suppose that $\mathcal{Y} = \{\eta_1, \dots, \eta_{|\mathcal{Y}|}\}$, where $\eta_1 > \dots > \eta_{|\mathcal{Y}|}$. Assume $D(r)$ and k grow unboundedly and $k/n = o(1)$. For any n , let $\tilde{D}(r) = D(r)$ if $r \leq k$ and $\tilde{D}(r) = 0$ otherwise. Assume $f(X) \in [a, b]$; $f(X)$ has a probability density function such that $\mathbb{P}(f(X) = s) > 0$ for all $s \in [a, b]$; $\Pr(Y = \eta_j) > 0$ and $\Pr(Y = \eta_j | f(X) = s)$ is a continuous function of s for all j . Then

$$\text{NDCG}_{\tilde{D}}(f, S_n) \xrightarrow{p} \frac{1}{\eta_1} \cdot \mathbb{E}[Y | \tilde{f}(X) = 1].$$

The last two theorems are for top- k , where $k/n = c$. We consider both logarithm discount and polynomial discount separately.

Theorem 12. Suppose that $\mathcal{Y} = \{\eta_1, \dots, \eta_{|\mathcal{Y}|}\}$, where $\eta_1 > \dots > \eta_{|\mathcal{Y}|}$. Let $k/n = c$ for some constant $c > 0$. Let $D(r) = \frac{1}{\log(1+r)}$. For any n , let $\tilde{D}(r) = D(r)$ if $r \leq k$ and $\tilde{D}(r) = 0$ otherwise. Assume $f(X) \in [a, b]$; $f(X)$ has a probability density function such that $\mathbb{P}(f(X) = s) > 0$ for all $s \in [a, b]$; $\Pr(Y = \eta_j) > 0$ and $\Pr(Y = \eta_j | f(X) = s)$ is a continuous function of s for all j . Then

$$\text{NDCG}_{\tilde{D}}(f, S_n) \xrightarrow{p} \frac{c \cdot \mathbb{E}[Y | \tilde{f}(X) \geq 1 - c]}{\sum_{j=1}^t \eta_j (R_j - R_{j-1}) + \eta_{t+1} (c - R_t)}.$$

where $R_0 = 0$; $R_j = \mathbb{P}(Y \geq \eta_j)$; t is defined by $R_t < c \leq R_{t+1}$.

Theorem 13. Let $D(r) = r^{-\beta}$ with $\beta \in (0, 1)$, and $\tilde{D}(r) = D(r)$ if $r \leq k$ and $\tilde{D}(r) = 0$ otherwise. Using the same notions and under the same conditions as in Theorem 12

$$\text{NDCG}_{\tilde{D}}(f, S_n) \xrightarrow{p} \frac{(1 - \beta) \int_{1-c}^1 \mathbb{E}[Y | \tilde{f}(X) = s] (1 - s)^{-\beta} ds}{\sum_{j=1}^t \eta_j (R_j^{1-\beta} - R_{j-1}^{1-\beta}) + \eta_{t+1} (c^{1-\beta} - R_t^{1-\beta})}.$$

4 Experimental Results

All theoretical results in this paper are proved under the assumption that the objects to rank are i.i.d. data. Often in real applications the data are not strictly i.i.d or even not random. Here we conduct experiments on a real dataset — Web search data. The aim is to see to what extent the behavior of the ranking measures on real datasets agree with our theory obtained under the i.i.d. assumption.

The dataset we use contains click-through log data of a search engine. We collected the clicked documents for 40 popular queries as test set, which are regarded as 40 independent ranking tasks. In each task, there are 5000 Web documents with clicks. To avoid heavy work of human labeling, we simply label each document by its click number according to the following rule. We assign relevancy $y = 2$ to documents with more than 1000 clicks, 1 to those with 100 to 1000 clicks, and 0 to the rest. In each task, we extracted 40 features for each item representing its relevance to the given query. A detail is how to construct S_n . In our theoretical analysis we assume S_n contains i.i.d. data. Since the goal of the experiments is to see how our theory works for real applications, we construct S_n as follows. For each query, there are totally 5000 documents which we denote by x_1, \dots, x_{5000} . Assume each document has a generating time. Without loss of generality we assume x_1 was generated earliest and x_{5000} latest. We set $S_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ for each $1 \leq n \leq 5000$. Such a construction simulates that in reality there may be increasing number of documents needed to rank by a search engine over time. We use three ranking functions in the experiments: a trained RankSVM model [20], a trained ListNet model [10], and a function chosen randomly. To be concrete, the random function is constructed as follows. For each $x \in \mathcal{X}$, we set $f(x)$ by choosing a number uniformly random from $[-1, 1]$. For the trained models (i.e., listNet and RankSVM), parameters are learned from a separate large training set construct in the same manner as the test set. Clearly, ListNet and RankSVM are relatively good ranking functions and the random function is bad.

We analyze the following typical NDCG type ranking measures by experiments:

- Standard NDCG: $D(r) = \frac{1}{\log(1+r)}$. See Figure 1, Theorem 1 and Theorem 2.
- NDCG with a feasible discount function: $D(r) = r^{-1/2}$. See Figure 2, Theorem 3 and Theorem 4.
- NDCG with too fast decay: $D(r) = 2^{-r}$. See Figure 3 and Theorem 6.

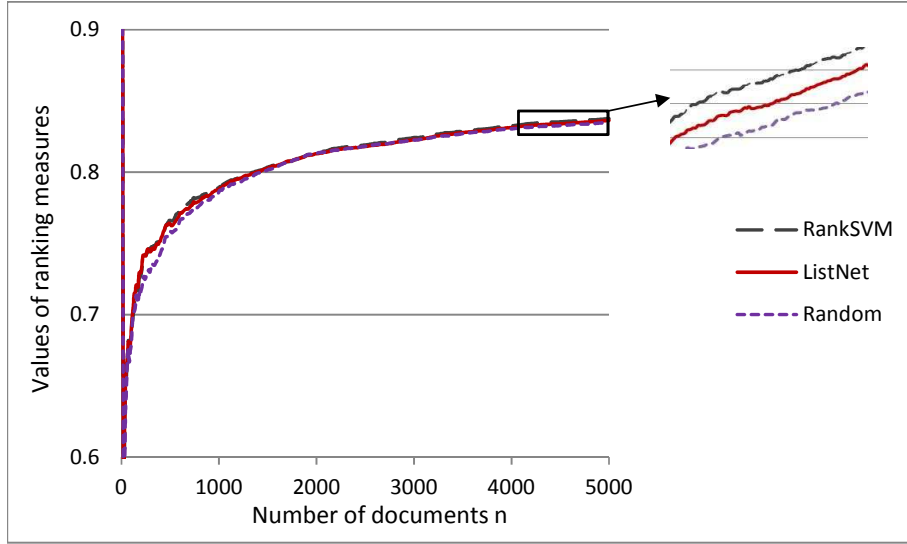


Figure 1: Standard NDCG: Converges to the same limit but distinguishes well the ranking functions.

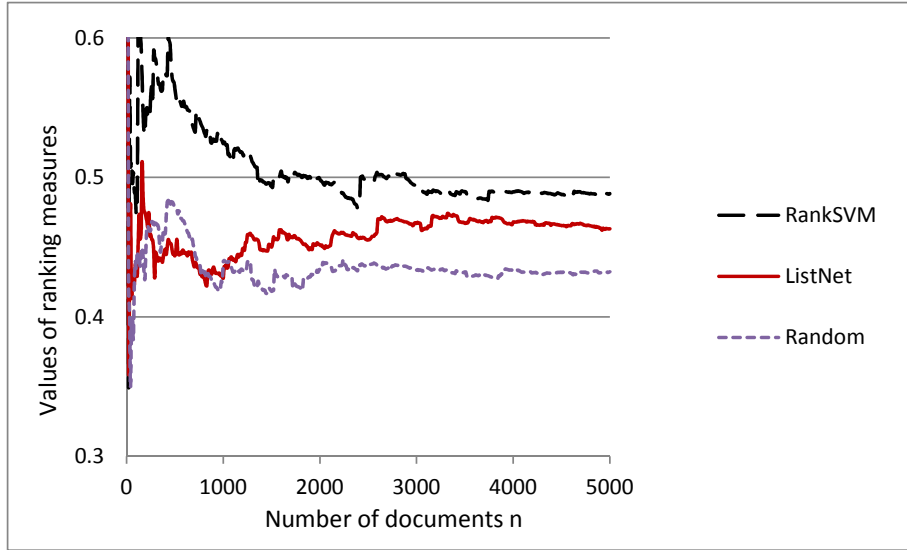


Figure 2: NDCG with feasible discount $D(r) = r^{-1/2}$: converges to different limits and distinguishes well the ranking functions.

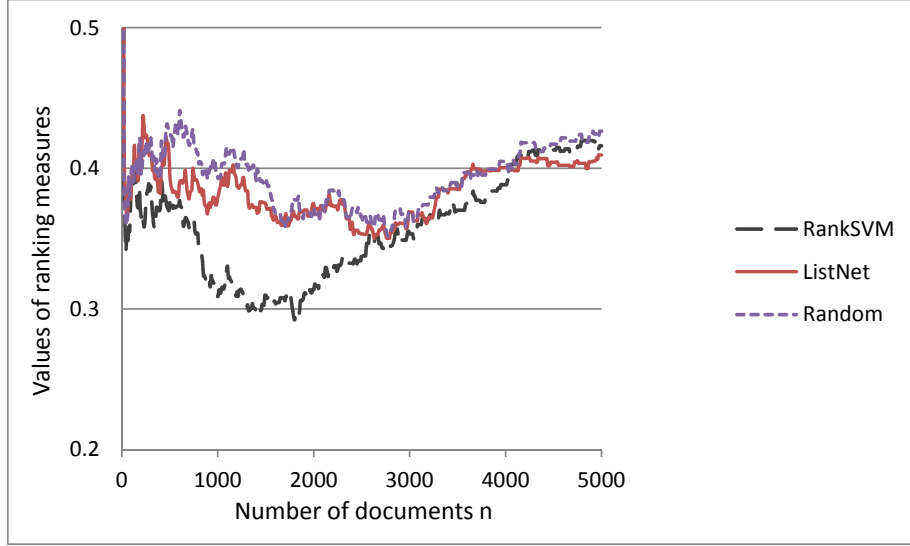


Figure 3: NDCG with too fast decay $D(r) = 2^{-r}$: does not converge; does not have good distinguishability power either.

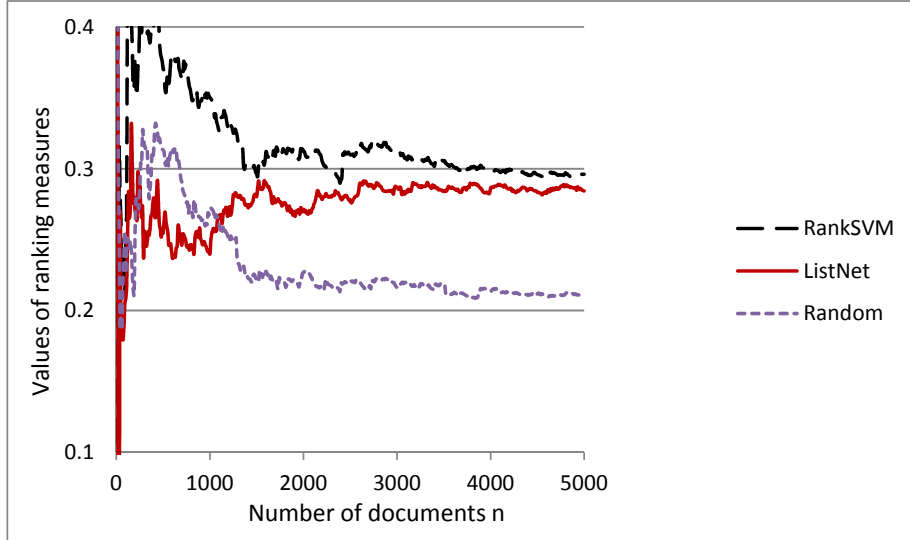


Figure 4: NDCG@ k ($D(r) = \frac{1}{\log(1+r)}$, $k = n/5$): distinguishes well the ranking functions.

- NDCG@k: $k = n/5$; $D(r) = \frac{1}{\log(1+r)}$. See Figure 4 and Theorem 8.

Figure 1 agrees well with Theorem 1 and Theorem 2. On the one hand, the NDCG measures of the three ranking functions are very close and seem to converge to the same limit. On the other hand, one can see from the enlarged part (we enlarge and stretch the vertical axis) in the figure that in fact the measures distinguish well the ranking functions.

Figure 2 demonstrates the result of NDCG with the feasible discount $r^{-1/2}$. In this experiment, it seems that the ranking measures of the three ranking functions converge to different limits and therefore distinguish them very well. In our experimental setting, it is not easy to find two ranking functions whose NDCG measures converge to the same limit. If one can find such a pair of ranking functions, it would be interesting to see how well the measure distinguish them.

Figure 3 shows the behavior of NDCG with a smooth discount which decays too fast. The measure cannot distinguish the three ranking functions very well. Even the randomly chosen function has an NDCG score similar to those of RankSVM and ListNet. From the figure, it is also likely that the measures do not converge.

Figure 4 depicts the result of NDCG@k, where k is a constant proportion of n . Before describing the result, let us first comparing Theorem 8 and Theorem 1. Note that although the discount are both the logarithmic one, NDCG@k for $k = cn$ can converge to different limits for different ranking functions, while standard NDCG always converges to 1. Figure 4 clearly demonstrate this result.

Acknowledgement

Liwei Wang would like to thank Kai Fan and Ziteng Wang for long and helpful discussions.

References

- [1] S. Agarwal, T. Graepel, R. Herbrich, S. Har-Peled, and D. Roth. Generalization bounds for the area under an ROC curve. 2004.
- [2] A. Al-Maskari, M. Sanderson, and P. Clough. The relationship between IR effectiveness measures and user satisfaction. In *SIGIR*, pages 773–774, 2007.
- [3] J. Aslam, E. Yilmaz, and V. Pavlu. The maximum entropy method for analyzing retrieval measures. In *Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 27–34. ACM, 2005.
- [4] R. Baeza-Yates and B. Ribeiro-Neto. *Modern Information Retrieval*, volume 82. Addison-Wesley New York, 1999.
- [5] M. Balcan, N. Bansal, A. Beygelzimer, D. Coppersmith, J. Langford, and G. Sorkin. Robust reductions from ranking to classification. *Machine learning*, 72(1):139–153, 2008.
- [6] P. Bartlett, M. Jordan, and J. McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- [7] D. Buffoni, C. Calauzenes, P. Gallinari, and N. Usunier. Learning scoring functions with order-preserving losses and standardized supervision. In *Proceedings of the 28th International Conference on Machine Learning*, 2011.
- [8] C. Burges, R. Ragno, and Q. V. Le. Learning to rank with nonsmooth cost functions. In *Advances in Neural Information Processing Systems 19: Proceedings of the 2006 Conference*, volume 19, page 193. The MIT Press, 2007.

- [9] C. Calauzènes, N. Usunier, and P. Gallinari. On the (non-)existence of convex, calibrated surrogate losses for ranking. In P. Bartlett, F. Pereira, C. Burges, L. Bottou, and K. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 197–205. 2012.
- [10] Z. Cao, T. Qin, T. Liu, M. Tsai, and H. Li. Learning to rank: from pairwise approach to listwise approach. In *ICML*, pages 129–136, 2007.
- [11] O. Chapelle, D. Metzler, Y. Zhang, and P. Grinspan. Expected reciprocal rank for graded relevance. In *Proceedings of the 18th ACM conference on Information and knowledge management*, pages 621–630. ACM, 2009.
- [12] S. J. Cléménçon and N. Vayatis. Empirical performance maximization for linear rank statistics. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems 21*, pages 305–312. 2009.
- [13] S. Cléménçon, G. Lugosi, and N. Vayatis. Ranking and empirical minimization of U-statistics. *The Annals of Statistics*, 36(2):844–874, 2008.
- [14] D. Cossock and T. Zhang. Statistical analysis of bayes optimal subset ranking. *Information Theory, IEEE Transactions on*, 54(11):5140–5154, 2008.
- [15] K. Crammer and Y. Singer. Pranking with ranking. In *Advances in Neural Information Processing Systems*, 2002.
- [16] W. Croft, D. Metzler, and T. Strohman. *Search engines: Information retrieval in practice*. Addison-Wesley, 2010.
- [17] J. Duchi, L. Mackey, and M. Jordan. On the consistency of ranking algorithms. In *Proceedings of the 27th International Conference on Machine Learning*, pages 327–334, 2010.
- [18] Y. Freund, R. Iyer, R. Schapire, and Y. Singer. An efficient boosting algorithm for combining preferences. *The Journal of Machine Learning Research*, 4:933–969, 2003.
- [19] J. Hájek, Z. Šidák, and P. Sen. *Theory of rank tests*. Academic press New York, 1967.
- [20] R. Herbrich, T. Graepel, and K. Obermayer. Large margin rank boundaries for ordinal regression. *Advances in Neural Information Processing Systems*, pages 115–132, 1999.
- [21] K. Järvelin and J. Kekäläinen. IR evaluation methods for retrieving highly relevant documents. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval*, pages 41–48. ACM, 2000.
- [22] K. Järvelin and J. Kekäläinen. Cumulated gain-based evaluation of IR techniques. *ACM Transactions on Information Systems (TOIS)*, 20(4):422–446, 2002.
- [23] T. Joachims. Optimizing search engines using clickthrough data. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 133–142. ACM, 2002.
- [24] E. Kanoulas and J. A. Aslam. Empirical justification of the gain and discount function for NDCG. In *Proceedings of the 18th ACM conference on Information and knowledge management*, pages 611–620. ACM, 2009.
- [25] M. Kendall. A new measure of rank correlation. *Biometrika*, 30(1/2):81–93, 1938.
- [26] R. Nallapati. Discriminative models for information retrieval. In *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 64–71. ACM, 2004.

- [27] P. Ravikumar, A. Tewari, and E. Yang. On NDCG consistency of listwise ranking methods. In *Proceedings of 14th International Conference on Artificial Intelligence and Statistics, AISTATS*, 2011.
- [28] C. Rudin. The p-norm push: A simple convex ranking algorithm that concentrates at the top of the list. *The Journal of Machine Learning Research*, 10:2233–2271, 2009.
- [29] T. Sakai. Evaluating evaluation metrics based on the bootstrap. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 525–532. ACM, 2006.
- [30] G. Sansone. *Orthogonal Functions*. Interscience Publishers Inc., New York, 1959.
- [31] A. Tewari and P. Bartlett. On the consistency of multiclass classification methods. *Journal of Machine Learning Research*, 8:1007–1025, 2007.
- [32] A. Turpin and F. Scholer. User performance versus precision measures for simple search tasks. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 11–18. ACM, 2006.
- [33] H. Valizadegan, R. Jin, R. Zhang, and J. Mao. Learning to rank by optimizing NDCG measure. *Advances in Neural Information Processing Systems*, 22:1883–1891, 2009.
- [34] E. Voorhees. Evaluation by highly relevant documents. In *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 74–82. ACM, 2001.
- [35] F. Xia, T. Liu, J. Wang, W. Zhang, and H. Li. Listwise approach to learning to rank: theory and algorithm. In *Proceedings of the 25th international conference on Machine learning*, pages 1192–1199. ACM, 2008.
- [36] Y. Yue, T. Finley, F. Radlinski, and T. Joachims. A support vector method for optimizing average precision. In *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 271–278. ACM, 2007.
- [37] T. Zhang. Statistical analysis of some multi-category large margin classification methods. *The Journal of Machine Learning Research*, 5:1225–1251, 2004.
- [38] T. Zhang. Statistical behavior and consistency of classification methods based on convex risk minimization. *Annals of Statistics*, pages 56–85, 2004.

A Proof of Theorem 2: the Key Lemmas

In this section we will prove Theorem 2. In fact we will prove a more complete result. The proof relies on a few key lemmas. In this section we only state these lemmas. Their proofs will be given in Appendix B. First we give a weaker definition of distinguishability, which guarantees that the ranking measure \mathcal{M} gives consistent comparison results for two ranking functions only in expectation.

Definition 4. Fix an underlying distribution P_{XY} . A pair of ranking functions f_0, f_1 is said to be distinguishable in expectation by a ranking measure \mathcal{M} , if there exist $b \in \{0, 1\}$ and a positive integer N such that for all $n \geq N$,

$$\mathbb{E}[\mathcal{M}(f_b, S_n)] > \mathbb{E}[\mathcal{M}(f_{1-b}, S_n)],$$

where the expectation is over the random draw of S_n .

Now we state a theorem which contains Theorem 2.

Theorem 14. Assume that $p = \Pr(Y = 1) > 0$. For every pair of ranking functions f_0, f_1 , Let $\bar{y}^{f_i}(s) = \Pr[Y = 1 | \tilde{f}_i(X) = s]$, $i = 0, 1$. Unless $\bar{y}^{f_0}(s) = \bar{y}^{f_1}(s)$ almost surely on $[0, 1]$, f_0, f_1 are distinguishable in expectation by standard NDCG whose discount is $D(r) = \frac{1}{\log(1+r)}$.

Moreover, if $\bar{y}^{f_0}(s)$ and $\bar{y}^{f_1}(s)$ are Hölder continuous in s , then unless $\bar{y}^{f_0}(s) = \bar{y}^{f_1}(s)$ almost everywhere on $[0, 1]$, f_0 and f_1 are consistently distinguishable by standard NDCG.

To prove Theorem 14, we need some notations.

Definition 5. Suppose $\mathcal{Y} = \{0, 1\}$. Let $\bar{y}^f(s) = \Pr[Y = 1 | \tilde{f}(X) = s]$. Also let $F(t) = \int_1^t D(s)ds$. We define the unnormalized pseudo-expectation $\tilde{N}_D^f(n)$ as

$$\tilde{N}_D^f(n) = \int_1^n \bar{y}^f(1 - s/n) D(s) ds = n \int_{\frac{1}{n}}^1 \bar{y}^f(1 - s) D(ns) ds.$$

Assume that $p = \Pr(Y = 1) > 0$. Define the normalized pseudo-expectation $N_D^f(n)$ as

$$N_D^f(n) = \frac{\tilde{N}_D^f(n)}{F(np)}.$$

The proof of the first part of Theorem 14 (i.e., distinguishable in expectation) relies on the following two key lemmas, whose proofs will be given in Appendix B.

Lemma 2. Let $D(r) = \frac{1}{\log(1+r)}$. Assume that $p = \Pr(Y = 1) > 0$. Then for every ranking function f ,

$$\left| \mathbb{E}[\text{NDCG}_D(f, S_n)] - N_D^f(n) \right| \leq \tilde{O}\left(n^{-1/3}\right).$$

Lemma 3. Let $D(r) = \frac{1}{\log(1+r)}$. Assume that $p = \Pr(Y = 1) > 0$. let $\bar{y}^{f_i}(s) = \Pr[Y = 1 | \tilde{f}_i(X) = s]$, $i = 0, 1$. Unless $\bar{y}^{f_0}(\cdot) = \bar{y}^{f_1}(\cdot)$ almost everywhere on $[0, 1]$, there must exist a nonnegative integer K and a constant $a \neq 0$, such that

$$\left| N_D^{f_0}(n) - N_D^{f_1}(n) - \frac{a}{\log^K n} \right| \leq O\left(\frac{1}{\log^{K+1} n}\right).$$

Lemma 2 says that the difference between the expectation of the NDCG measure of a ranking function and its pseudo-expectation is relatively small; while Lemma 3 says that the difference between the pseudo-expectations of two essentially different ranking functions are much larger.

To prove the “moreover” part of Theorem 14 (i.e., consistently distinguishable), we need the following key lemma, whose proof will be given in Section B. The lemma states that with high probability the NDCG measure of a ranking function is very close to its pseudo-expectation.

Lemma 4. Let $D(r) = \frac{1}{\log(1+r)}$. Assume that $p = \Pr(Y = 1) > 0$. Suppose the ranking function f satisfies that $\bar{y}^f(s) = \Pr(Y = 1 | \tilde{f}(X) = s)$ is Hölder continuous with constants $\alpha > 0$ and $C > 0$. That is, $|\bar{y}^f(s) - \bar{y}^f(s')| \leq C|s - s'|^\alpha$ for all $s, s' \in [0, 1]$. Then

$$\Pr\left[\left|\text{NDCG}_D(f, S_n) - N_D^f(n)\right| \geq 5Cp^{-1}n^{-\min(\alpha/3, 1)}\right] \leq O\left(e^{-n^{1/4}}\right).$$

Proof. of Theorem 14 That f_0 and f_1 are strictly distinguishable in expectation by standard NDCG is straightforward from Lemma 2 and Lemma 3. That f_0 and f_1 are strictly distinguishable with high probability follows immediately from Lemma 4, Lemma 3 and the observation that $\sum_{n \geq N} e^{-n^{1/4}} \leq O\left(N^{3/4}e^{-N^{1/4}}\right) \leq O\left(e^{-N^{1/5}}\right)$. \square

B Proofs of the Key Lemmas in Appendix A

In this section, we give proofs of the three key lemmas in Appendix A (i.e., Lemma 2, Lemma 3 and Lemma 4) used to prove Theorem 2 and Theorem 14.

To prove the key lemmas, we need a few technical claims, whose proofs will be given in Appendix C. We first give four claims that will be used in the proof of Lemma 2.

Claim 1. For any $s \in [0, 1]$,

$$\sum_{r=1}^n \mathbb{P}[\tilde{f}(x_{(r)}^f) = s] = n. \quad (8)$$

Claim 2. Recall that the DCG ranking measure with respect to discount $D(\cdot)$ was defined as

$$\text{DCG}_D(f, S_n) = \sum_{r=1}^n y_{(r)}^f D(r). \quad (9)$$

Let $D(r) = \frac{1}{\log(1+r)}$, and $\bar{y}^f(s) = \Pr[Y = 1 | \tilde{f}(X) = s]$. Then

$$\mathbb{E}[\text{DCG}_D(f, S_n)] = \sum_{r=1}^n \frac{1}{\log(1+r)} \int_0^1 \mathbb{P}[\tilde{f}(x_{(r)}^f) = 1-s] \bar{y}^f(1-s) ds. \quad (10)$$

Claim 3. For any positive integer n , define $E_{n,r} = [\frac{r}{n} - n^{-1/3}, \frac{r}{n} + n^{-1/3}]$ ($r \in [n]$). Then for any $r \in [n]$,

$$\Pr[1 - \tilde{f}(x_{(r)}^f) \in E_{n,r}] \geq 1 - 2e^{-n^{1/3}}. \quad (11)$$

Claim 4. Let $\mathcal{Y} = \{0, 1\}$. Assume $D(r) = \frac{1}{\log(1+r)}$. Let $F(t) = \int_1^t D(s) ds$. Assume also $p = \Pr[Y = 1] > 0$. Then for every sufficiently large n , with probability $(1 - 2e^{-2n^{1/3}})$ the following inequality holds.

$$\left| \text{NDCG}_D(f, S_n) - \frac{\text{DCG}_D(f, S_n)}{F(np)} \right| \leq O(n^{-1/3}). \quad (12)$$

Now we are ready to prove Lemma 2.

Proof. of Lemma 2. By the definition of $\tilde{N}_D^f(n)$ (see Definition 5) and eq.(8), we have

$$\tilde{N}_D^f(n) = n \int_{\frac{1}{n}}^1 \frac{\bar{y}^f(1-s) ds}{\log(1+ns)} = \sum_{r=1}^n \int_{\frac{1}{n}}^1 \frac{\bar{y}^f(1-s) \mathbb{P}[\tilde{f}(x_{(r)}^f) = 1-s]}{\log(1+ns)} ds. \quad (13)$$

By eq. (10) in Claim 2 and eq.(13), and note that $\bar{y}^f(s) \leq 1$, we obtain

$$\begin{aligned} & \left| \mathbb{E}[\text{DCG}_D(f, S_n)] - \tilde{N}_D^f(n) \right| \\ & \leq \sum_{r=1}^n \left| \int_{\frac{1}{n}}^1 \bar{y}^f(1-s) \mathbb{P}[\tilde{f}(x_{(r)}^f) = 1-s] \left(\frac{1}{\log(1+r)} - \frac{1}{\log(1+ns)} \right) ds \right| + \frac{1}{n} \sum_{r=1}^n \frac{1}{\log(1+r)} \\ & \leq \sum_{r=1}^n \left| \int_{[\frac{1}{n}, 1] \setminus E_{n,r}} \mathbb{P}[\tilde{f}(x_{(r)}^f) = 1-s] \left| \frac{1}{\log(1+r)} - \frac{1}{\log(1+ns)} \right| ds \right| \\ & + \sum_{r=1}^n \int_{E_{n,r} \cap [\frac{1}{n}, 1]} \left| \frac{1}{\log(1+r)} - \frac{1}{\log(1+ns)} \right| ds + O\left(\frac{1}{\log n}\right). \end{aligned} \quad (14)$$

We next bound the two terms in the RHS of the last inequality of (14) separately. By Claim 3, the first term can be upper bounded by

$$2e^{-n^{1/3}} \sum_{r=1}^n \sup_{s \in [\frac{1}{n}, 1] \setminus E_{n,r}} \left| \frac{1}{\log(1+r)} - \frac{1}{\log(1+ns)} \right| \leq \frac{2}{\log 2} n e^{-2n^{1/3}}. \quad (15)$$

For the second term in the RHS of the last inequality of (14), it is easy to check that the following two inequalities hold:

$$\forall r > n^{2/3}, \quad \sup_{s \in E_{n,r} \cap [\frac{1}{n}, 1]} \left| \frac{1}{\log(1+r)} - \frac{1}{\log(1+ns)} \right| \leq \frac{n^{2/3}}{(1+r) \log^2(1+r)} + o\left(\frac{n^{2/3}}{(1+r) \log^2(1+r)}\right). \quad (16)$$

$$\forall r \leq n^{2/3}, \quad \sup_{s \in E_{n,r} \cap [\frac{1}{n}, 1]} \left| \frac{1}{\log(1+r)} - \frac{1}{\log(1+ns)} \right| \leq \frac{1}{\log 2}. \quad (17)$$

Combining (14), (15), (16) and (17), we obtain

$$\left| \mathbb{E}[\text{DCG}_D(f, S_n)] - \tilde{N}_D^f(n) \right| \leq \frac{2ne^{-2ne^{1/3}}}{\log 2} + O\left(\frac{n^{2/3}}{\log 2} + \sum_{r=n^{2/3}}^n \frac{n^{2/3}}{(1+r) \log^2(1+r)}\right) \leq \tilde{O}(n^{2/3}). \quad (18)$$

Finally, observe that $F(np) = \text{Li}(1+np)$, where Li is the offset logarithmic integral function. By Claim 4 and the well-known fact $\text{Li}(n) \sim \frac{n}{\log n}$, we have the following inequality and this completes the proof.

$$\left| \mathbb{E}[\text{NDCG}_D(f, S_n)] - \frac{\mathbb{E}[\text{DCG}_D(f, S_n)]}{\text{Li}(1+np)} \right| \leq \tilde{O}(n^{-1/3}) + O(e^{-2n^{1/3}}). \quad (19)$$

□

We next turn to prove Lemma 3. We need the following three claims.

Claim 5. For sufficiently large n ,

$$\int_0^{\frac{2}{n}} \log^k x dx = O\left(\frac{\log^k n}{n}\right). \quad (20)$$

Claim 6. Fix an integer $k \in \mathbb{N}^* = \{0\} \cup \mathbb{N}$. For sufficiently large n ,

$$\int_{\frac{2}{n}}^1 \frac{|\log^k x| dx}{(\log(nx))^{k+1}} \leq O\left(\frac{1}{\log^{k+1} n}\right). \quad (21)$$

Claim 7. $\text{span}\left(\{\log^k x\}_{k \geq 0}\right)$, is dense in $L^2[0, 1]$.

Now we are ready to prove Lemma 3.

Proof. of Lemma 3. Let $\Delta y(s) = \bar{y}^{f_0}(s) - \bar{y}^{f_1}(s)$. By the definition of normalized pseudo expectation (see definition 5) and the fact that $|\Delta y(s)| \leq 1$, we have

$$\begin{aligned} N_D^{f_0}(n) - N_D^{f_1}(n) &= \frac{n}{\text{Li}(1+np)} \int_{\frac{1}{n}}^1 \frac{\Delta y(1-s) ds}{\log(1+ns)} \\ &= \frac{n}{\text{Li}(1+np)} \int_{\frac{2}{n}}^1 \frac{\Delta y(1-s) ds}{\log(1+ns)} + O\left(\frac{1}{\text{Li}(n)}\right). \end{aligned} \quad (22)$$

Expanding $\frac{1}{\log(1+ns)}$ at the point ns , we obtain

$$\left| \int_{\frac{2}{n}}^1 \frac{\Delta y(1-s)ds}{\log(1+ns)} - \int_{\frac{2}{n}}^1 \frac{\Delta y(1-s)ds}{\log n + \log s} \right| \leq \int_{\frac{2}{n}}^1 \frac{ds}{ns \log^2(ns)} \leq O\left(\frac{\log n}{n}\right). \quad (23)$$

Expanding $\frac{1}{\log n + \log s}$ at point $\log n$, we have that for all $m \in \mathbb{N}^*$, the following holds:

$$\begin{aligned} & \left| \int_{\frac{2}{n}}^1 \frac{\Delta y(1-s)ds}{\log n + \log s} - \sum_{j=1}^m \frac{(-1)^{j-1}}{\log^j n} \int_{\frac{2}{n}}^1 \Delta y(1-s) \log^{j-1} s \, ds \right| \\ &= \left| \int_{\frac{2}{n}}^1 \frac{\Delta y(1-s) \log^m s \, ds}{(\log n + \xi_{n,s})^{m+1}} \right| \leq \int_{\frac{2}{n}}^1 \frac{|\Delta y(1-s) \log^m s| \, ds}{(\log n + \log s)^{m+1}} \leq O\left(\frac{1}{\log^{m+1} n}\right). \end{aligned} \quad (24)$$

Note in above derivation that $\xi_{n,s} \in (\log s, 0)$, and the last inequality is due to Claim 6.

Furthermore, by Claim 7, unless $\Delta y(s) = 0$ a.e., there exist constants $k \in \mathbb{N}^*$ and $a \neq 0$ such that

$$(-1)^k \int_0^1 \Delta y(1-s) \log^k s \, ds = a. \quad (25)$$

Let K be the smallest integer k that Eq. (25) holds. Combining (22), (23), (24), and (25) and noting Claim 5, we have the following and this completes the proof.

$$\left| N_D^{f_0}(n) - N_D^{f_1}(n) - \frac{a}{\log^K n} \right| \leq O\left(\frac{\log^K n}{n}\right) + O\left(\frac{1}{\log^{K+1} n}\right).$$

□

To prove the last key lemma, we need the following claim.

Claim 8. Let $D(r) = \frac{1}{\log(1+r)}$. Let $F(t) = \int_1^t D(r)dr$. Assume $\bar{y}^f(s)$ is Hölder continuous with constants α and C . Then

$$\left| \sum_{r=1}^n \bar{y}^f(1-r/n) D_r - \tilde{N}_D^f(n) \right| \leq Cn^{-\alpha/3} F(n) + D(1) + |D'(1)|. \quad (26)$$

Now we prove the last key lemma.

Proof. of Lemma 4. Let x_1, \dots, x_n be instances i.i.d. drawn according to P_X . Let $\tilde{x}_{(r)} = \tilde{f}(x_{(r)}^f)$ and by definition $\tilde{x}_{(1)} \geq \tilde{x}_{(2)} \geq \dots \geq \tilde{x}_{(n)}$. By Chernoff bound, for every r with probability $2e^{-2n^{1/3}}$ we have $|\tilde{x}_{(r)} - (1-r/n)| > n^{-1/3}$. A union bound over r then yields

$$\Pr \left[\forall r \in [n], \left| \tilde{x}_{(r)} - \left(1 - \frac{r}{n}\right) \right| \leq n^{-1/3} \right] \geq 1 - 2ne^{-2n^{1/3}}. \quad (27)$$

Since y^f is Hölder continuous with constants α and C , eq. (27) implies

$$\Pr \left[\left| \sum_{r=1}^n \bar{y}^f(\tilde{x}_{(r)}) D(r) - \sum_{r=1}^n \bar{y}^f(1-r/n) D(r) \right| \leq Cn^{-\alpha/3} \cdot \sum_{r=1}^n D(r) \right] \geq 1 - 2ne^{-2n^{1/3}}. \quad (28)$$

Combining Claim 8 and eq. (28), and note that $|D'(1)| + D(1) \leq 10$ we have

$$\Pr \left[\left| \sum_{r=1}^n \bar{y}^f(\tilde{x}_{(r)}) D(r) - \tilde{N}_D^f(n) \right| \leq 2Cn^{-\alpha/3} \cdot F(n) + 10 \right] \geq 1 - 2ne^{-2n^{1/3}}. \quad (29)$$

Fix x_1, \dots, x_n . Let $x_{(1)}^f, \dots, x_{(n)}^f$ be the induced ordered sequence. Also let $\tilde{x}_{(r)} = \tilde{f}(x_{(r)}^f)$. Recall that $\bar{y}^f(s) = \mathbb{E}[Y | \tilde{f}(X) = s]$. Thus $\sum_{r=1}^n \bar{y}^f(\tilde{x}_{(r)}) D(r)$ is the expectation of $\text{DCG}_D(f, S_n) = \sum_{r=1}^n y_{(r)}^f D(r)$ conditioned on the fixed values $\tilde{x}_{(1)}, \dots, \tilde{x}_{(n)}$. Also observe that conditioning on $\tilde{x}_{(1)}, \dots, \tilde{x}_{(n)}$, $y_{(r)}^f$ ($r = 1, \dots, n$) are independent. By Hoeffding's inequality and taking into consideration that x_1, \dots, x_n are arbitrary and $(D(r))^2 \leq D(r)$ for all r , we have for every $\epsilon > 0$

$$\Pr \left[\left| \text{DCG}_D(f, S_n) - \sum_{r=1}^n \bar{y}^f(\tilde{x}_{(r)}) D(r) \right| \geq \epsilon \right] \leq 2 \exp \left(-\frac{2\epsilon^2}{F(n)} \right). \quad (30)$$

Set $\epsilon = F(n)^{2/3}$ in eq. (30) and combine eq. (29), we have

$$\Pr \left[\left| \text{DCG}_D(f, S_n) - \tilde{N}_D^f(n) \right| > 2Cn^{-\alpha/3} F(n) + 2F(n)^{2/3} \right] \leq 2ne^{-2n^{1/3}} + 2e^{-2F(n)^{1/3}}. \quad (31)$$

Simple calculations yields

$$\Pr \left[\left| \frac{\text{DCG}_D(f, S_n)}{F(np)} - N_D^f(n) \right| > 4Cp^{-1} n^{-\min(\alpha/3, 1)} \right] \leq 2ne^{-2n^{1/3}} + 2e^{-2F(n)^{1/3}}. \quad (32)$$

Combining eq. (12) and (32) The lemma follows. \square

C Proof of the Technical Claims in Appendix B

Here we give proofs of the technical claims by which we prove the three key lemmas in Section B.

Proof. of Claim 1.

Recall that for each $i \in [n]$, $\tilde{f}(x_i)$ is uniformly distributed on $[0, 1]$; and $x_{(1)}^f, \dots, x_{(n)}^f$ are just reordering of x_1, \dots, x_n . Thus

$$\sum_{r=1}^n \mathbb{P}[\tilde{f}(x_{(r)}^f) = s] = \sum_{i=1}^n \mathbb{P}[\tilde{f}(x_i) = s] = n.$$

\square

Proof. of Claim 2.

We have

$$\begin{aligned} \mathbb{E}[\text{DCG}_D(f, S_n)] &= \sum_{r=1}^n D(r) \mathbb{E}[y_{(r)}^f] \\ &= \sum_{r=1}^n \frac{1}{\log(1+r)} \mathbb{E} \left[\mathbb{E}[y_{(r)}^f | \tilde{f}(x_{(r)}^f)] \right] \\ &= \sum_{r=1}^n \frac{1}{\log(1+r)} \int_0^1 \mathbb{P}[\tilde{f}(x_{(r)}^f) = s] \bar{y}^f(s) ds. \end{aligned} \quad (33)$$

\square

Proof. of Claim 3. Just observe that $\tilde{f}(x_{(r)}^f)$ is the r -th order statistic (r -th largest) of n uniformly distributed random variables on $[0, 1]$. Chernoff bound yields the result. \square

Proof. of Claim 4.

Let $l = \sum_{(x,y) \in S_n} \mathbb{I}[y = 1]$ be the number of $y = 1$ in S_n . Since S_n is sampled i.i.d. and $\Pr[Y = 1] = p$, by Chernoff bound we have

$$\Pr \left[\left| l/n - p \right| > n^{-1/3} \right] \leq 2e^{-2n^{1/3}}. \quad (34)$$

Thus with probability at least $1 - 2e^{-2n^{1/3}}$

$$\begin{aligned} & \left| \text{NDCG}_D(f, S_n) - \frac{\text{DCG}_D(f, S_n)}{F(np)} \right| \\ &= \left| \frac{\text{DCG}_D(f, S_n)}{l} - \frac{\text{DCG}_D(f, S_n)}{F(np)} \right| \\ &\leq \text{DCG}_D(f, S_n) \cdot \max \left(\left| \frac{1}{F(n(p - n^{-1/3}))} - \frac{1}{F(np)} \right|, \left| \frac{1}{F(n(p + n^{-1/3}))} - \frac{1}{F(np)} \right| \right). \end{aligned}$$

Recall that $F(t) = \int_1^t \frac{1}{\log(1+r)} dr$, $p > 0$; and observe that $\text{DCG}_D(f, S_n) \leq F(n)$. Taylor expansion of $\frac{1}{F((p \pm n^{-1/3})n)}$ at np and some simple calculations yields the result. \square

Proof. of Claim 5.

Integration by part we have,

$$\int \log^k x dx = k! \sum_{j=0}^k (-1)^{k-j} \frac{x \log^j x}{j!} + C. \quad (35)$$

The claim follows. \square

Proof. of Claim 6.

Changing variable by letting $x = n^{-t}$ we have

$$\begin{aligned} & \int_{\frac{2}{n}}^1 \frac{|\log^k x| dx}{(\log(nx))^{k+1}} \\ &= \int_0^{1 - \frac{\log 2}{\log n}} \frac{t^k}{(1-t)^{k+1}} e^{-t \log n} dt \\ &= \int_0^{1/2} \frac{t^k}{(1-t)^{k+1}} e^{-t \log n} dt + \int_{1/2}^{1 - \frac{\log 2}{\log n}} \frac{t^k}{(1-t)^{k+1}} e^{-t \log n} dt. \end{aligned} \quad (36)$$

Now we upper bound the two terms in the last line of eq. (36) separately. For the first term we have

$$\begin{aligned} & \int_0^{1/2} \frac{t^k}{(1-t)^{k+1}} e^{-t \log n} dt \leq 2^{k+1} \int_0^{1/2} t^k e^{-t \log n} dt \\ &\leq \frac{2^{k+1}}{(\log n)^{k+1}} \int_0^\infty \tau^k e^{-\tau} d\tau \leq \frac{2^{k+1} \Gamma(k+1)}{(\log n)^{k+1}} \\ &= O\left(\frac{1}{(\log n)^{k+1}}\right), \end{aligned} \quad (37)$$

where Γ is the gamma function, and the last inequality is due to that k is a fixed integer.

For the second term we have

$$\begin{aligned} & \int_{1/2}^{1-\frac{\log 2}{\log n}} \frac{t^k}{(1-t)^{k+1}} e^{-t \log n} dt \leq \left(\frac{\log n}{\log 2} \right)^{k+1} \int_{1/2}^1 e^{-t \log n} dt \\ & \leq \frac{1}{2} \cdot \frac{1}{\sqrt{n}} \cdot \left(\frac{\log n}{\log 2} \right)^{k+1} = \tilde{O}\left(\frac{1}{\sqrt{n}}\right), \end{aligned} \quad (38)$$

where in \tilde{O} we hide the $\text{polylog}(n)$ terms.

Combining (37) and (38) we complete the proof. \square

Proof. of Claim 7.

We only need to show that for any $f \in L^2[0, 1]$, if

$$\int_0^1 f(x) \log^k x dx = 0, \quad k = 0, 1, \dots \quad (39)$$

then $f = 0$ a.e. on $[0, 1]$.

Let $t = -\log x$, then eq.(39) becomes

$$\int_0^\infty f(e^{-t}) t^k e^{-t} dt = 0, \quad k = 0, 1, \dots$$

Note that Laguerre polynomials form a complete basis of $L^2[0, \infty)$ (cf. [30], p.349), thus $\{t^k\}_{k \geq 0}$ is complete in $L^2[0, \infty)$ with respect to measure e^{-t} . The claim follows. \square

Proof. of Claim 8.

$$\begin{aligned} & \left| \sum_{r=1}^n \bar{y}^f(1-r/n) D(r) - \tilde{N}_D^f(n) \right| \\ &= \left| \sum_{r=1}^n \bar{y}^f(1-r/n) D(r) - \int_1^n \bar{y}^f(1-s/n) D(s) ds \right| \\ &= \left| \sum_{r=1}^{n-1} \int_r^{r+1} (\bar{y}^f(1-r/n) D(r) - \bar{y}^f(1-s/n) D(s)) ds \right| + \bar{y}^f(0) D(n) \\ &\leq \left| \sum_{r=1}^{n-1} \int_r^{r+1} \bar{y}^f(1-s/n) (D(r) - D(s)) ds \right| \\ &\quad + \sum_{r=1}^{n-1} \int_r^{r+1} \left| \bar{y}^f(1-r/n) - \bar{y}^f(1-s/n) \right| D(r) ds + \bar{y}^f(0) D(n) \\ &\leq \sum_{r=1}^{n-1} \int_r^{r+1} |D(r) - D(s)| ds + Cn^{-\alpha/3} \sum_{r=1}^{n-1} D(r) + D(n) \\ &\leq \sum_{r=1}^{n-1} |D'(r)| + Cn^{-\alpha/3} F(n) + D(n) \\ &\leq Cn^{-\alpha/3} F(n) + |D'(1)| + \sum_{r=2}^n |D'(r)| + D(n) \\ &\leq Cn^{-\alpha/3} F(n) + |D'(1)| + D(1) - D(n) + D(n) \\ &= Cn^{-\alpha/3} F(n) + |D'(1)| + D(1). \end{aligned}$$

Note that the sixth and the seventh line are both because $|D'(r)|$ is monotone decreasing; and second line from bottom is because $D(r)$ is monotone decreasing. \square

D Proof of the Convergence Theorems

In this section we give the proof of the theorems considering convergence of NDCG with various discount and cut-off.

First we give the proof of Theorem 1, i.e., the standard NDCG converges to 1 almost surely for every ranking function.

Proof. of Theorem 1. For notational simplicity we only prove for the case $\mathcal{Y} = \{0, 1\}$. Generalization is straightforward. Recall that $S_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ consists of n i.i.d. instance-label pairs drawn from an underlying distribution P_{XY} . Let $p = \Pr(Y = 1)$. Also let $l = \sum_{i=1}^n y_i$. If $p = 0$, the theorem trivially holds. Suppose $p > 0$, by Chernoff bound we have

$$\Pr\left(\left|\frac{l}{n} - p\right| > n^{-1/3}\right) \leq 2e^{-2n^{1/3}}.$$

For fixed n , conditioned on the event that $\left|\frac{l}{n} - p\right| \leq n^{-1/3}$, by the definition of NDCG, it is easy to see that

$$\begin{aligned} \text{NDCG}_D(f, S_n) &= \frac{\sum_{r=1}^n y_{(r)}^f \frac{1}{\log(1+r)}}{\sum_{r=1}^l \frac{1}{\log(1+r)}} \\ &\geq \frac{\sum_{r=n-l+1}^n \frac{1}{\log(1+r)}}{\sum_{r=1}^l \frac{1}{\log(1+r)}} \\ &\geq \frac{\text{Li}(n+1) - \text{Li}(n(1-p + n^{-1/3}) + 1)}{\text{Li}(n(p + n^{-1/3}) + 1)} - o(1) \\ &\geq 1 - o(1), \end{aligned} \tag{40}$$

where $\text{Li}(t) = \int_2^t \frac{d\tau}{\log \tau}$ is the offset logarithmic integral function; and the last step in eq.(40) is due to the well-known fact that $\text{Li}(t) \sim \frac{t}{\log t}$. Thus for any $\epsilon > 0$, and for any sufficiently large n , conditioned on the event that $\left|\frac{1}{n} \sum_{i=1}^n y_i - p\right| \leq n^{-1/3}$, we have

$$\text{NDCG}_D(f, S_n) \geq 1 - \epsilon.$$

Also recall that $\text{NDCG}_D(f, S_n) \leq 1$. We have, for any $\epsilon > 0$ and every sufficiently large n

$$\Pr(|\text{NDCG}_D(f, S_n) - 1| \geq \epsilon) \leq 2e^{-2n^{1/3}}.$$

Since $\sum_{n \geq 1} 2e^{-2n^{1/3}} < \infty$, by Borel-Cantelli lemma $\text{NDCG}_D(f, S_n)$ converges to 1 almost surely. \square

Next we give details of the other feasible discount functions as well as the cut-off versions. In particular, we provide proofs of Theorems 3, 5, 7, 8, 9. The proofs of these five theorems are quite similar. We only prove Theorem 3 to illustrate the ideas. The proof of the other four theorems require only minor modifications.

The proof of Theorem 3 relies on the following lemma, which is similar to Lemma 4.

Lemma 5. Let $D(r) = r^{-\beta}$ for some $\beta \in (0, 1)$. Assume that $p = \Pr(Y = 1) > 0$. If the ranking function f satisfies that $\bar{y}^f(s) = \Pr(Y = 1 | f(X) = s)$ is continuous, then for every $\epsilon > 0$ the following inequality holds for all sufficiently large n :

$$\Pr \left[\left| \text{NDCG}_D(f, S_n) - N_D^f(n) \right| \geq 5p^{-1}\epsilon \right] \leq o(1).$$

Proof. of Theorem 3. The theorem follows from Lemma 5 and simple calculations of $\lim_{n \rightarrow \infty} N_D^f(n)$. We omit the details. \square

Proof. of Lemma 5. The proof is simple modification of the proof of Lemma 4. Note that the difference of Lemma 5 from Lemma 4 is that here we do not assume y^f is Hölder continuous. We only assume it is continuous.

Next observe that Claim 8 holds for $D(r) = r^{-\beta}$ ($0 < \beta < 1$) as well. Because in the proof of Claim 8, we only use two properties of $D(r)$. That is, $D(r)$ is monotone decreasing and $|D'(r)|$ is monotone decreasing. Clearly $D(r) = r^{-\beta}$ satisfies these properties. But here $\bar{y}^f(s)$ is merely continuous rather than Hölder continuous. Thus we have a modified version of Claim 8. That is, for every $\epsilon > 0$, the following holds for all sufficiently large n :

$$\left| \sum_{r=1}^n \bar{y}^f(1 - r/n) D(r) - \tilde{N}_D^f(n) \right| \leq \epsilon F(n) + D(1) + |D'(1)|.$$

The rest of the proof are almost identical to Lemma 4. We omit the details. \square

Finally, we give the proof of Theorem 6, i.e., if the discount decays substantially faster than r^{-1} , then the NDCG measure does not converge. Moreover, every pair of ranking functions are not strictly distinguishable with high probability by the measure.

Proof. of Theorem 6. For notational simplicity we give a proof for $|\mathcal{Y}| = 2$ and $\mathcal{Y} = \{0, 1\}$. It is straightforward to generalize it to other cases.

In fact, we only need to show that for every ranking function f , there are constants $a, b, c > 0$ with $a > b$, such that for all sufficiently large n ,

$$\Pr[\text{NDCG}_D(f, S_n) \geq a] \geq c$$

and

$$\Pr[\text{NDCG}_D(f, S_n) \leq b] \geq c$$

both hold. Once we prove this, by definition the ranking measure does not converge (in probability). Also, it is clear that for every pair of ranking functions, there is at least a constant probability that the ranking measure of the two functions are “overlap”. Therefore distinguishability is not possible.

For sufficiently large n , fix any x_1, \dots, x_n . According to the assumption, the probability that the top-ranked m data all have label 1 is at least $(\delta/2)^m$, where m is the minimal integer such that

$$\sum_{r=1}^m D(r) \geq \frac{2}{3} \sum_{r=1}^{\infty} D(r).$$

Clearly we have

$$\Pr \left(\text{NDCG}_D(f, S_n) \geq \frac{2}{3} \mid x_1, \dots, x_n \right) \geq (\delta/2)^m.$$

On the other hand, the probability that the top-ranked m elements all have label 0 and there are at least m elements in the list that have label 1 is at least $(\delta/2)^{2m}$. Note that

$$\frac{\sum_{r=m+1}^n D(r)}{\sum_{r=1}^m D(r)} \leq \frac{1}{2}.$$

Thus we have

$$\Pr[\text{NDCG}_D(f, S_n) \leq \frac{1}{2} \mid x_1, \dots, x_n] \geq (\delta/2)^{2m}.$$

Since x_1, \dots, x_n are arbitrary, the theorem follows. \square

E Proof of Distinguishability for NDCG with $r^{-\beta}$ ($\beta \in (0, 1)$) Discount

Here we give the proof of Theorem 4, i.e., NDCG with $r^{-\beta}$ ($0 < \beta < 1$) discount has the power of distinguishability.

Proof. of Theorem 4. The proof of distinguishability for polynomial discount is much easier than that of the logarithmic discount, because in the former case the pseudo-expectation has very simple form. If f_0 and f_1 satisfy the first condition $\int_0^1 \Delta y(s)(1-s)^{-\beta} ds \neq 0$, then the theorem is trivially true since $\text{NDCG}(f_0, S_n)$ and $\text{NDCG}(f_1, S_n)$ converge to different limits. So we only need to prove the theorem assuming that $\int_0^1 \Delta y(s)(1-s)^{-\beta} ds = 0$ and the second condition holds. The proof is similar to Theorem 2 by using the pseudo-expectation. We have the next two lemmas for discount $D(r) = r^{-\beta}$, $\beta \in (0, 1)$.

Lemma 6. *Let $D(r) = r^{-\beta}$, $\beta \in (0, 1)$. Suppose that $\bar{y}^{f_0}(s)$ and $\bar{y}^{f_1}(s)$ are continuous. Also assume that $\int_0^1 \Delta y(s)(1-s)^{-\beta} ds = 0$ and $\Delta y(1) \neq 0$. Then we have*

$$\left| N_D^{f_0}(n) - N_D^{f_1}(n) \right| \geq \left| \frac{\Delta y(1)}{2p^{1-\beta}} \right| \cdot n^{-(1-\beta)}. \quad (41)$$

Proof.

$$\begin{aligned} N_D^{f_0}(n) - N_D^{f_1}(n) &= \frac{n}{F(np)} \int_{1/n}^1 \Delta y(1-s) \cdot (ns)^{-\beta} ds \\ &= \frac{1-\beta}{p^{1-\beta}} \int_{1/n}^1 \Delta y(1-s) \cdot s^{-\beta} ds \\ &= -\frac{1-\beta}{p^{1-\beta}} \int_0^{1/n} \Delta y(1-s) \cdot s^{-\beta} ds. \end{aligned}$$

Since Δy is continuous, for any $\delta > 0$ there exists $\epsilon > 0$ such that for all $x \in [1-\epsilon, 1]$, $|\Delta y(x) - \Delta y(1)| \leq \delta$. Consequently, for sufficiently large n ,

$$\left| \int_0^{1/n} \Delta y(1-s) \cdot s^{-\beta} ds - \Delta y(1) \cdot \int_0^{1/n} s^{-\beta} ds \right| \leq \delta \cdot \int_0^{1/n} s^{-\beta} ds.$$

Let $\delta = \Delta y(1)/2$, we then have

$$\left| \int_0^{1/n} \Delta y(1-s) \cdot s^{-\beta} ds - \frac{\Delta y(1)}{1-\beta} \cdot n^{-(1-\beta)} \right| \leq \frac{\Delta y(1)}{2(1-\beta)} \cdot n^{-(1-\beta)}.$$

The lemma follows. \square

Lemma 7. *Let $D(r) = r^{-\beta}$, $\beta \in (0, 1)$. Assume that $p = \Pr(Y = 1) > 0$. If the ranking function f satisfies that $\bar{y}^f(s) = \Pr(Y = 1 \mid f(X) = s)$ is Hölder continuous with constants $\alpha > 0$ and $C > 0$ That is, $|\bar{y}^f(s) - \bar{y}^f(s')| \leq C|s - s'|^\alpha$ for all $s, s' \in [0, 1]$. Then*

$$\Pr \left[\left| \text{NDCG}_D(f, S_n) - N_D^f(n) \right| \geq 5Cp^{-1}n^{-\min(\alpha/3, 1)} \right] \leq O \left(e^{-n^{(1-\beta)/3}} \right).$$

Proof. The proof is almost the same as the proof of Lemma 4

□

The theorem follows immediately from Lemma 6 and Lemma 7.

□